

Low-dimensional Query Projection based on Divergence Minimization Feedback Model for Ad-hoc Retrieval

Javid Dadashkarimi

Masoud Jalili Sabet

Heshaam Faili

Azadeh Shakery

School of ECE, College of Engineering, University of Tehran, Iran
{dadashkarimi,jalili.masoud,hfaili, shakery}@ut.ac.ir

ABSTRACT

Low-dimensional word vectors have long been used in a wide range of applications in natural language processing. In this paper we shed light on estimating query vectors in ad-hoc retrieval where a limited information is available in the original query. Pseudo-relevance feedback (PRF) is a well-known technique for updating query language models and expanding the queries with a number of relevant terms. We formulate the query updating in low-dimensional spaces first with rotating the query vector and then with scaling. These consequential steps are embedded in a query-specific projection matrix capturing both angle and scaling. In this paper we propose a new but not the most effective technique necessarily for PRF in language modeling, based on the query projection algorithm. We learn an embedded coefficient matrix for each query, whose aim is to improve the vector representation of the query by transforming it to a more reliable space, and then update the query language model. The proposed embedded coefficient divergence minimization model (ECDMM) takes top-ranked documents retrieved by the query and obtains a couple of positive and negative sample sets; these samples are used for learning the coefficient matrix which will be used for projecting the query vector and updating the query language model using a softmax function. Experimental results on several TREC and CLEF data sets in several languages demonstrate effectiveness of ECDMM. The experimental results reveal that the new formulation for the query works as well as state-of-the-art PRF techniques and outperforms state-of-the-art PRF techniques in a TREC collection in terms of MAP@5, and P@10 significantly.

Keywords

low-dimensional word vectors, query embedding, query projection, feedback model, information retrieval.

1. INTRODUCTION

Low-dimensional word vectors have long been tailored as useful resources for understanding machine-readable texts. latent semantic analysis (LSI), probabilistic latent semantic analysis (PLSI),

non-negative matrix factorization (NMF), and latent Dirichlet allocation (LDA) are statistical techniques to this end [7, 13, 31, 3].

LSI reduces high-dimensional count vectors to low-dimensional latent semantic space [7]; PLSI provides a statistical formulation on LSA [13]; NMF is tailored over a term-document weighting matrix for term recommendation and query re-weighting [31]; and LDA assumes Dirichlet allocation for both document-topic and term-topic distributions and introduce a better generative model to this end [3]. Although their effectiveness in low-dimensional text modeling, all the methods rely on the simplified bag of word assumption; Recently, low-dimensional word vectors constructed from a feed-forward neural networks with single hidden layer have been used successfully for text modeling in a wide range of applications [22, 5, 12]. The neural network-based vectors captures positions of the words in a context and indirectly embed multiple n-grams in the final model [22]. The constructed vectors have been used in document classification [16], term-by-term expansion [1], sentiment analysis [26], named entity recognition (NER) [28], semantic role labeling [5], and query re-weighting [33]. A linear neural network learns a number of hidden variables for a word, capturing both semantic and syntactic information [22, 33]. A further convolutional layer is used for estimating vectors for pieces of texts, sentences, documents, and queries [5]. Nevertheless, simple *max*, *min*, and *avg* aggregation layers are used instead successfully [29, 18].

Word embedding has limited applications in ad-hoc retrieval but, it is likely to have a burgeon in this area during the next years [21]. However, it has been shown that query-document vector similarity, in which there is not document length normalization nor term frequency consideration, degrades the retrieval performance and thus it is recommend to use these vectors within state-of-the-art retrieval frameworks [33, 29, 6]. Estimating robust language models for documents and queries is a required task in this area. Although it has been proposed a number of interesting works in low-dimensional text modeling [8, 19], but as far as we know, there is no in-detail study investigating on low-dimensional query language modeling where a few number of keywords being posed by the users.

In this paper we shed light on estimating query vectors in ad-hoc retrieval where a limited information about a user intention is available in the original query. Therefore it is required to propose a novel technique for building/updating the low-dimensional query vectors which is the focus of this paper. We formulate the query modeling/updating in low-dimensional spaces first with rotating the query vector and then with scaling. These consequential steps can be embedded in a query-specific projection matrix capturing both angle and scaling. Thus we propose to learn this matrix for each query and then project the query vector to a more relevant low-

dimensional space. Learning this matrix requires low-dimensional vectors of a number of relevant/irrelevant vectors to the query. To this aim, first we extract a couple of word sets from top-ranked documents retrieved in response to the original query and then find a matrix projecting the query vector close to the relevant vectors and far away from the irrelevant ones. The projected query vector is then be used for building a query language model after applying a softmax/sigmoid function. The obtained query language model can be tailored in a statistical language modeling framework, the state-of-the-art retrieval framework.

Experimental results on several data sets from TREC and CLEF in English, French, Spanish, German, and Persian demonstrate the effectiveness of the proposed method. The experimental results reveal that the new formulation for the query vector modeling/updating works as well as state-of-the-art PRF techniques and even better in a few number of collections. The proposed method outperforms all competitive baselines in terms of MAP in a TREC collection and a German collection. The performances in the French collection and the Spanish one are very competitive to NMF.

2. PREVIOUS WORKS

2.1 Pseudo-relevance Feedback

Top-ranked documents $F = \{d_1, d_2, \dots, d_{|F|}\}$ in response to a query $q = \{q_1, q_2, \dots, q_m\}$ have long been used as helpful resources for expanding the query [17, 20]. Lavrenko et al., introduced relevance models for updating the query. The RM1 method models the query as:

$$p(w|\theta_q) \propto \sum_{d \in F} p(w|\theta_d) p(\theta_d) \prod_{i=1}^{|q|} p(q_i|\theta_d) \quad (1)$$

where θ_d is the language model of document $d \in F$. The RM2 models the query in another way as:

$$p(w|\theta_q) \propto p(w) \prod_{i=1}^{|q|} \sum_{d \in F} p(q_i|\theta_d) \frac{p(w|\theta_d) p(\theta_d)}{p(w)}. \quad (2)$$

The obtained relevance models can be interpolated with the original query as follows:

$$p(w|\theta'_q) = (1 - \alpha) p(w|\theta_q) + \alpha p_{\text{ml}}(w|q) \quad (3)$$

where $p_{\text{ml}}(w|q)$ is the maximum likelihood estimation of the original query. The interpolated model for RM1 and RM2 are known as RM3 and RM4 respectively [14].

Zhai & Lafferty in [32] introduced the mixture model (MIXTURE) based on an expectation maximization algorithm for modeling the feedback documents.

$$t^{(n)}(w) = \frac{(1 - \lambda) p_\lambda^{(n)}(w|\theta_F)}{(1 - \lambda) p_\lambda^{(n)}(w|\theta_F) + \lambda p(w|\mathcal{C})} \quad (4)$$

$$p_\lambda^{(n+1)}(w|\theta_F) = \frac{\sum_{j=1}^n c(w; d_j) t^{(n)}(w)}{\sum_i \sum_{j=1}^n c(w_i; d_j) t^{(n)}(w)} \quad (5)$$

where $t^{(n)}(w)$ is topicality of the word w at n -th iteration. $c(w_i; d_j)$ indicates the count of the word w_i in document d_j . The authors introduced another model for feedback modeling based on

divergence minimization in [32] as follows:

$$p(w|\hat{\theta}_F) \propto \exp \left(\frac{1}{1 - \lambda} \frac{1}{|F|} \sum_i \log p(w|\hat{\theta}_{d_i}) - \frac{\lambda}{1 - \lambda} \log p(w|\mathcal{C}) \right) \quad (6)$$

where λ is a controlling constant. Later Lv & Zhai modified this framework in [20] and introduced maximum entropy divergence minimization model (MEDMM) that aims to find $\hat{\theta}_F$ as follow:

$$\arg \min_{\theta} \sum_{d \in F} \alpha_d H(\theta_F, \theta_d) - \lambda H(\theta_F, \theta_C) - \beta H(\theta_F) \quad (7)$$

where $H(\theta_1, \theta_2)$ is the cross-entropy between θ_1 and θ_2 and $H(\theta)$ is the entropy of θ .

Zamani et al., [31] proposed a new technique for PRF based on non-negative matrix factorization. The proposed relevance feedback based on matrix factorization (RFMF) first builds $A \in \mathbb{R}_+^{m \times n}$ based on occurrences of terms in $d \in F$ where $m = |F| + 1$ and n is the number of unique words in F (the $|F| + 1$ -th row belongs to the original query). Second, RFMF aims to decompose $A^{m \times n}$ to $U \in \mathbb{R}_+^{m \times r}$ and $V \in \mathbb{R}_+^{r \times n}$. The final product of $U^{m \times r}$ and $V^{r \times n}$ re-weights the query language model.

2.2 Low-dimensional Vectors

Low-dimensional representations of words are tailored in a variety of tasks in natural language processing [5]. To learn these vectors a common approach is to predict the context of each word and then aim at minimizing a loss function. This method is known as skip-gram negative sampling and can be interpreted as a binary regression task as follows:

$$\arg \min_{\theta} \sum_{(w,c) \in D \cup D'} \log \left(\left(\frac{1}{1 + \exp^{-\mathbf{v}_w^T \mathbf{v}_c}} \right)^z \left(\frac{1}{1 + \exp^{\mathbf{v}_w^T \mathbf{v}_c}} \right)^{1-z} \right) \quad (8)$$

where $\mathbf{v}_w \in \mathbb{R}^{n \times 1}$ and $\mathbf{v}_c \in \mathbb{R}^{n \times 1}$ are the vectors of a word and its context respectively. z indicates if this sample (w, c) is a valid sample ($z = 1$) or not ($z = 0$) [9]. [23] introduced global word vector (GloVe) as follows:

$$\arg \min_{\theta} \sum_{i,j}^V f(X_{ij}) (v_{w_i}^T v_{w_j} + b_i + b_j - \log X_{ij}) \quad (9)$$

where $X_{ij} \in \mathbb{R}^{V \times V}$ is the co-occurrence matrix, b_i and b_j are constant biases, and $f(X)$ is defined as

$$f(x) = \begin{cases} (x/x_{\max})^\alpha & \text{if } x < x_{\max} \\ 1 & \text{other} \end{cases} \quad (10)$$

Estimating low-dimension vectors of sentences and documents is another interesting area in the litterateur. There are a couple of approaches to this end. The first approach is based on offline algorithms trying to estimate a document/sentence representation based on a number of available low-dimensional word vectors [29, 4]. Usually a convolutional layer is used for the final estimation [5]. The second approach is based on learning a representation vector for each document/sentence during learning the vectors of the words [19]. Indeed, the proposed paragraph-to-vector is suitable for pre-defined documents or queries. Herein, we do not provide this approach first for lower performance of direct query-document low-dimensional vector similarity in ad-hoc retrieval and second for the streaming essence of the queries in this area.

Kusner et al., in [16] exploited distributional document representations in a number of text classification tasks. ALMasri et

al., in [1] introduced a term-by-term expansion approach based on the low-dimensional vector similarities. Kiros et al., in [15] investigated on expanding a text by incorporating vector similarity of the candidates with averaged vector of the embedded words in the text. The obtained vectors are used in sentiment classification, cross-lingual document classification, blog authorship attribution, and conditional word similarity. However, there is a number of works with advanced neural networks [25, 26]. Socher et al., introduce recurrent neural networks (RNN) for sentence-formed queries/tweets.

[29] proposed an information retrieval framework based on the low-dimensional vectors. The authors demonstrated that the low-dimensional vectors are not yet effective in the vector-space ad-hoc retrieval frameworks where we compute document-query low-dimensional vector similarity. Zheng et al., incorporated the word vectors in a supervised technique for re-weighting terms in probabilistic language model and BM25 [33]. Grbovic et al., uses query logs for building query models. The obtained models are used for query prediction and advertisement [10]. A number of works investigate on using the embedded vectors in cross-lingual environments [6, 29, 2]. [6] employed an offline projection algorithm to bridge the gap between the languages. The authors incorporated the vector similarities for building a query language model. [29] uses an on-line shuffling approach to this aim. The low-dimensional vectors are learned on a large comparable corpora after shuffling the words of each alignment. Dadashkarimi et al., demonstrated that the language model obtained by the projected vectors from different languages outperforms the shuffling approach [6]. Bengio et al., introduced BIBOWA, a fast on-line technique for learning multilingual word vectors, that learns the vectors from parallel corpora instead of bilingual dictionaries [2].

3. EMBEDDED COEFFICIENTS FOR QUERY PROJECTION

In this section we propose a novel technique for query language modeling based on low-dimensional word vectors. The proposed ECDMM takes the low-dimensional query vector and a number of relevant/irrelevant embedded vectors for its rotation and scaling. To this end, we aim to find a coefficient matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ for projecting the query model \mathbf{v}_q to a more relevant space by minimizing f as follows:

$$\begin{aligned} f(\mathbf{W}) = & \sum_{w_n \in F} \frac{\alpha}{2} \|\mathbf{W}^T \mathbf{v}_q - \mathbf{v}_{w_n}\|^2 \\ & - \sum_{\bar{w}_n \in F} \frac{\lambda}{2} \|\mathbf{W}^T \mathbf{v}_q - \mathbf{v}_{\bar{w}_n}\|^2 \\ & - \frac{\beta}{2} \|\mathbf{W}^T \mathbf{W}\| \end{aligned} \quad (11)$$

where $\mathbf{v}_q \in \mathbb{R}^{n \times 1}$ is the query vector, $\mathbf{v}_{w_n} \in \mathbb{R}^{n \times 1}$ is the vector of a relevant sample from F , and $\mathbf{v}_{\bar{w}_n} \in \mathbb{R}^{n \times 1}$ is the vector of a non-relevant sample from F . α is a controlling constant for the positive samples and λ is for the negative ones. β is another constant parameter for the regularization term $\|\mathbf{W}^T \mathbf{W}\|$.

The query vector is built from averaging the vectors of query words as follows:

$$[\mathbf{v}_q]_j \leftarrow \frac{1}{m} \sum_{1 \leq i \leq m=|q|} [\mathbf{v}_{q_i}]_j \quad (12)$$

Equation 11 has similar components to MEDMM [20]. The first part captures the same essence as the cross-entropy between the feedback model θ_F and the positive sample model θ_w in MEDMM.

It tries to minimize the distance between the query model and the positive samples. The second part also captures the effect of the negative samples on θ_F . It tries to maximize the distance between the query model and the negative samples. And, $\mathbf{W}^T \mathbf{W}$ acts as a regularization term for \mathbf{W} in the model. Positive samples are drawn according to the following distribution:

$$w^+ \sim \frac{(1 - \lambda)p_{\text{ml}}(w|\theta_F)}{(1 - \lambda)p_{\text{ml}}(w|\theta_F) + \lambda p(w|\mathcal{C})} \quad (13)$$

where θ_F and θ_C are unigram feedback distribution and unigram collection distribution respectively. λ is set near to 0.9 empirically for penalizing common words in \mathcal{C} [32]. Negative samples are drawn as follows:

$$w^- \sim p(w|\theta_C; F)^{\frac{3}{4}} \quad (14)$$

which is the unigram language model of the feedback documents raised to the power of $\frac{3}{4}$. This power is used for increasing the chance of appearing rare words [9]. To find the optimum value of \mathbf{W} we use the stochastic gradient descent algorithm as follows:

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \eta \frac{\partial f}{\partial \mathbf{W}} \quad (15)$$

where η is a constant learning rate. Since $f(\mathbf{W})$ is a differentiable single-variable quadratic function $\frac{\partial f}{\partial \mathbf{W}}$ can be obtained as follows:

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{W}} = & \alpha \sum_{w_n \in F} (\mathbf{W}^T \mathbf{v}_q - \mathbf{v}_{w_n}) \mathbf{v}_q^T \\ & - \lambda \sum_{\bar{w}_n \in F} (\mathbf{W}^T \mathbf{v}_q - \mathbf{v}_{\bar{w}_n}) \mathbf{v}_q^T \\ & - \beta \mathbf{W} \end{aligned} \quad (16)$$

Finally, the query vector \mathbf{V}_q can be projected to the new space as follows:

$$\hat{\mathbf{v}}_q = \mathbf{W}^T \mathbf{v}_q \quad (17)$$

3.1 Building Embedded Query Model

In this subsection we shed light on a number of approaches for building a query language model. Experiments in [29] demonstrate that using direct query-document vector similarity degrades the retrieval performance in monolingual document retrieval. Therefore, we aim to use a state-of-the-art retrieval framework as a core of document scoring and the low-dimensional vectors for query language modeling.

We believe that $\hat{\mathbf{v}}_q$ represents a better semantic direction of the query and therefore, we aim to find a feedback model based on similarity of the projected query with the words from the feedback documents. Herein, we examine a couple of functions to compute this similarity. First we use the sigmoid function as follows:

$$\nabla_1(\hat{\mathbf{v}}_q, \mathbf{v}_{w_n}) = \frac{1}{1 + e^{-\hat{\mathbf{v}}_q^T \mathbf{v}_{w_n}}} \quad (18)$$

Second we can use the cosine similarity to this end:

$$\nabla_2(\hat{\mathbf{v}}_q, \mathbf{v}_{w_n}) = \frac{\hat{\mathbf{v}}_q^T \mathbf{v}_{w_n}}{\|\hat{\mathbf{v}}_q\| \|\mathbf{v}_{w_n}\|} \quad (19)$$

To build a feedback language model we can use both ∇_1 and ∇_2 within a softmax layer as follows:

$$p(w_n|\hat{\theta}_F) = \frac{e^{\nabla(\hat{\mathbf{v}}_q, \mathbf{v}_{w_n})}}{\sum_{w_n} e^{\nabla(\hat{\mathbf{v}}_q, \mathbf{v}_{w_n})}} \quad (20)$$

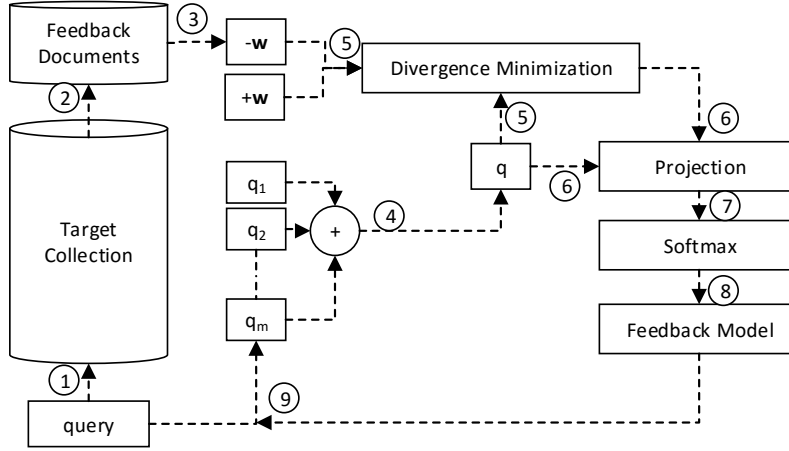


Figure 1: Whole process of low-dimensional query language modeling and feedback modeling.

Table 1: Collections Characteristics

ID	Lang.	Collection	Queries (title only)	#docs	#qrels
AP	English	Associated Press 88-89	TREC 1-3 Ad-Hoc Track, Q:51-200	164,597	15,838
ROB356	English	Los Angeles Times 1994, plus Glasgow Herald 1995	CLEF 2003-2004-2006, Q:141-350	170,153	4,327
SP	Spanish	EFE 1994	CLEF 2002, Q:91-140	215,738	2,854
DE	German	Frankfurter Rundschau 94, SDA 94, Der Spiegel 94-95	CLEF 2002-03, Q:91-140	225,371	1,938
FR	French	Le Monde 94, SDA French 94-95	CLEF 2002-03, Q:251-350	129,806	3,524
FA	Persian	Hamshahri 1996-2002	CLEF 2008-09, Q:551-650	166,774	9,625

Equation 20 captures only specificity of the feedback terms [27]. Term frequency is another useful metric for scoring the terms [27]. Herein, we use the following weighted softmax to incorporate both of these metrics:

$$p(w_n|\hat{\theta}_F) \propto \frac{a_{w_n} e^{\nabla(\hat{q}_q, \mathbf{v}_{w_n})}}{\sum_{w_n} a_{w_n} e^{\nabla(\hat{q}_q, \mathbf{v}_{w_n})}} \quad (21)$$

where $a_{w_n} = c(w_n, F)$ is the frequency of w in F . The obtained feedback model can be interpolated with the original query as shown in Eq. 3.

Figure 1 shows whole the process of query language modeling and feedback modeling.

4. EXPERIMENTS

4.1 Experimental Setup

Overview of the used collections are provided in Table 1. In all experiments, we use the language modeling framework with the KL-divergence retrieval model and Dirichlet smoothing with $\mu = 1000$. All European documents and queries are stemmed by the Porter stemmer and the Persian collection is remained intact [11, 24]. Therefore, we might have different vectors for cognate words. Stopwords are removed in all the experiments¹. The Lemur toolkit²

¹We use the stopwords lists and the normalizing techniques available at <http://www.unine.ch/info/clef/>.

²<http://www.lemurproject.org/>

Table 4: Comparing ECDMM with word2vec and ECDMM with GloVe. The * superscript shows a significant difference (2-tail t-test, $p \leq 0.05$).

	AP		ROB356	
	word2vec	GloVe	word2vec	GloVe
MAP	0.3330*	0.3245	0.3866	0.3895
P@5	0.4792	0.4738	0.4405	0.4405
P@10	0.4631	0.4510	0.3810	0.3856

is employed as the retrieval engine in our experiments. ECDMM is compared with the following methods: (1) maximum likelihood estimation of query (MLE): $p(w|\theta_q) = \frac{c(w, q)}{|q|}$ where $c(w, q)$ is the count of term w in the query; (2) RM3, (3) RM4, and (4) MEDMM explained in Section 2.1.

α in Equation 3 is set via 2-fold cross validation over topics of each collection and number of blind relevant documents is assumed $|F|=10$. All free parameters $\alpha, \lambda, \beta, n^+$, and n^- are fixed for all experiments after learning on a small sub-set of topics from the AP collection. Empirically we fixed the parameters to $\alpha = 0.8$, $\lambda = 0.05$, $\beta = 0.01$, $n^+ = 40$, and $n^- = 100$ in all the experiments. \mathbf{W} in Equation 16 is initialized with random values in $[-1, 1]$; η is set to a small value which also decreases after each iteration. The iterations terminate when the changes are very small or the number of iterations meets 1000.

The words' vectors computed with word2vec introduced in [22];

Table 2: Investigating ECDMM performance using different vector similarity methods. * indicates the weighted softmax function introduced in Eq. 21 (see also Eq. 18 and Eq. 19 for more details).

ID	AP			ROB356		
	MAP	P@5	P@10	MAP	P@5	P@10
SIGMOID/ ∇_1	0.3136	0.4698	0.4470	0.3434	0.4039	0.3667
COSIN/ ∇_2	0.3150	0.4711	0.4483	0.3486	0.4183	0.3699
SIGMOID*/ ∇_1	0.3241	0.4711	0.4544	0.3474	0.4078	0.3706
COSIN*/ ∇_2	0.3330	0.4792	0.4631	0.3866	0.4405	0.3810

Table 3: Comparison of different feedback methods. Superscripts 1/2/3/4/5 indicate that the MAP improvements are statistically significant compared to MLE/MIXTURE/RM3/RM4/MEDMM respectively (2-tail t-test, $p \leq 0.05$). The bold values in each column show the highest performance in terms of the corresponding metric for each collection.

ID	AP			ROB356			DE		
	MAP	P@5	P@10	MAP	P@5	P@10	MAP	P@5	P@10
MLE	0.2643	0.451	0.4262	0.3721	0.4366	0.3719	0.348	0.532	0.458
MIXTURE	0.3106	0.4450	0.4232	0.3781	0.4366	0.3758	0.4334	0.536	0.488
RM3	0.3187	0.4470	0.4294	0.4037	0.4405	0.3902	0.4346	0.5560	0.494
RM4	0.2875	0.4208	0.3876	0.3789	0.4392	0.3752	0.3632	0.5400	0.4720
MEDMM	0.3269	0.4551	0.4289	0.3908	0.4523	0.3876	0.3878	0.5400	0.4980
RFMF	0.3296	0.4577	0.4356	0.4096	0.451	0.402	0.4248	0.524	0.488
ECDMM	0.3330 ¹²³⁴⁵	0.4792	0.4631	0.3866 ¹²⁴	0.4405	0.3810	0.4369 ¹⁴⁵	0.552	0.4960
ID	FA			FR			SP		
	MAP	P@5	P@10	MAP	P@5	P@10	MAP	P@5	P@10
MLE	0.3554	0.584	0.561	0.3936	0.5212	0.4556	0.488	0.664	0.578
MIXTURE	0.3934	0.606	0.577	0.4119	0.499	0.4616	0.5161	0.644	0.598
RM3	0.4036	0.618	0.594	0.4115	0.4889	0.4556	0.5388	0.6520	0.6120
RM4	0.3721	0.596	0.586	0.4031	0.5172	0.4606	0.5053	0.6800	0.594
MEDMM	0.3585	0.59	0.56	0.4125	0.5253	0.4707	0.5268	0.688	0.608
RFMF	0.392	0.6100	0.5940	0.4219	0.5051	0.4626	0.5459	0.6640	0.6180
ECDMM	0.3950 ¹⁴⁵	0.6040	0.5730	0.4217 ¹²³⁴⁵	0.5152	0.4717	0.5384 ¹²⁴⁵	0.652	0.602

size of the window, number of negative samples, and size of the vectors are set to typical values of 10, 45, and 100 respectively. Vectors of GloVe are extracted from Wikipedia and Gigawords 5 with 6 billion tokens and 400k words³. The number of dimensions in GloVe is set 100 similarly.

4.2 Comparing word2vec and GloVe

In this section we aim to choose a standard low-dimensional vectors for our primary experiments. Therefore, we investigate performance of the proposed method using both word2vec and GloVe [22, 23]. Query vectors are built using the vectors obtained from word2vec and GloVe for the query terms in AP and ROB356. After that, the same configuration is used for evaluating the performance of ECDMM based on these vectors.

As shown in Table 4 ECDMM with word2vec outperforms ECDMM with GloVe in AP significantly. Differences in ROB356 are not statistically significant and thus we opted word2vec vectors in the rest of experiments.

4.3 Performance Comparison and Discussion

Table 2 shows experimental results on different vector similarity methods. According to the results, cosine similarity with a weighted softmax function works better than other similarity metrics. Indeed, the experimental results demonstrate that considering only the vector similarity captures the IDF value of the terms. But, incorporating this with counts of the terms works more better. Therefore we opt cosine similarity with the weighted softmax

function in the rest of experiments for comparison and parameter sensitivity.

All the experimental results are provided in Table 3. As shown in the table, the proposed ECDMM outperforms a few number of baselines in terms of MAP and works as well as the state-of-the-art feedback language models. The results show 26.0%, 3.8%, 25.5%, 11.1%, 7.1%, 10.3% improvements in AP, ROB356, DE, FA, FR, and SP respectively in terms of MAP compared to MLE. In AP and DE, ECDMM outperforms all the baselines in terms of all the metrics. Generally, the results in P@5 and P@10 are very competitive and there is not any specific method that outperforms all other ones in all the collections.

RFMF outperforms all the baselines and the proposed ECDMM in SP in all metrics. Differences in this collection are statistically significant. On the other hand, ECDMM outperforms RFMF in AP and DE significantly. Differences in other collections are not statistically significant.

As stated before, we have not used any stemmer for Persian and thus we might have different vectors for cognate words. But, the results belonging to FA demonstrate that the obtained vectors are suitable enough for our model.

As discussed in Section 3, ECDMM takes advantage of the first step of MIXTURE for positive sampling and the idea of MEDMM for divergence minimization. The experimental results reveal that ECDMM is more effective than these methods and captures both topicality and entropy. However, RM3 is a strong baseline and works better than the proposed method in ROB356, FA, and SP in terms of MAP. It is noteworthy that $p(\theta_d)$ in RM1/RM3 (see Eq. 1)

³<http://nlp.stanford.edu/projects/glove/>

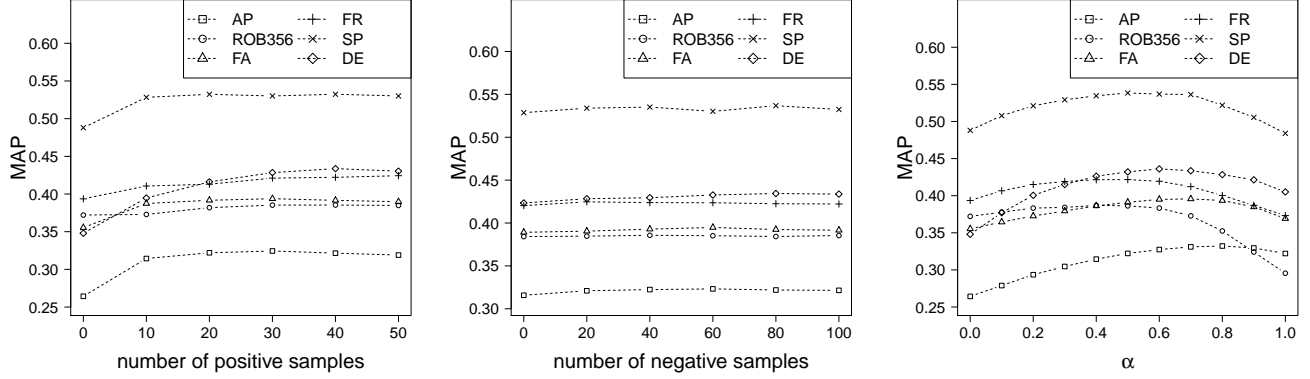


Figure 2: MAP sensitivity of ECDMM to n^+ , n^- , and α (feedback coefficient).

Table 5: Query language models for "airbus subsidy" by different feedback models. Terms are stemmed by the Porter stemmer.

MIXTURE		RM3		RM4		MEDMM		RFMF		ECDMM	
airbu	0.0558	airbu	0.0317	said	0.0297	airbu	0.1121	airbu	0.0299	airbu	0.1001
subsidi	0.0274	said	0.0197	airbu	0.0246	subsidi	0.0668	subsidi	0.0280	subsidi	0.0610
us	0.0227	airbu	0.0169	subsidi	0.0179	said	0.0381	consortium	0.0226	industri	0.0203
govern	0.0219	govern	0.0155	s	0.01769	us	0.0367	manufactur	0.0199	aircraft	0.0173
trade	0.0167	us	0.0153	us	0.01397	govern	0.0326	aircraft	0.0170	econom	0.0169
aircraft	0.0160	s	0.0117	govern	0.0137	consortium	0.0322	germani	0.0170	european	0.0158
west	0.0158	west	0.0110	aircraft	0.0096	aircraft	0.0271	blohm	0.0176	yeutter	0.0140
industri	0.0156	trade	0.0106	trade	0.0093	trade	0.0255	boelkow	0.0176	germani	0.0138
european	0.0132	industri	0.0095	will	0.0059	west	0.0088	spain	0.0160	daimler	0.0137
consortium	0.0126	aircraft	0.0094	aircraft	0.0052	state	0.0085	yeutter	0.0160	consortium	0.0126

plays key role in its efficiency and it would be interesting to study this effect on a_{w_n} in Eq. 21. Nevertheless, the differences are not statistically significant.

Table 5 shows top-10 stemmed expansion terms and the weights of obtained query model by the methods for 'airbus subsidy' from topics of AP. It shows that MIXTURE and ECDMM are more successful in purifying the feedback model from common words (see *s* and *said* in the lists). However, terms like *us* never appeared in ECDMM although it is semantically related to the query; As shown in the table, MEDMM and ECDMM weight the original query more than others (see *airbu* and *subsidi* in the lists). Therefore, higher value of α works well for the proposed model in the AP collection (see Fig.2 and [20]).

4.4 Parameter Sensitivity

In this section we investigate the sensitivity of the proposed method to the number of positive and negative samples and the feedback coefficient (see Eq. 3). To this aim, one parameter is fixed to its optimum value and different values are tested for the other one. As shown in Figure 2 both parameters n_+ and n_- work stable in all the collections. However as shown in the figure, the performance of the system is less reliable to n_- than n_+ . The reason might be due to existing a variety of topics in the non-relevant low-dimensional space. Nevertheless, in difficult queries, which is not the focus of the current work, it is necessary to consider negative feedback as well. In this kind of circumstances, F contains fewer number of relevant documents and a retrieval system is required to use this negative information for query modification [30].

Vulic et al., have shown that the retrieval performance is not sensitive to the number of dimensions considerably [29]. Therefore,

herein we fixed this parameter to the typical value of 100 and investigate sensitivity of other parameters.

5. CONCLUSION AND FUTURE WORKS

In this paper, we propose a query language model using low-dimensional query projection. We use a set of positive and negative samples from the top-ranked documents, retrieved by the query, to learn an embedded coefficient matrix. The query vector, which got transformed by the coefficient matrix, is then used to expand the original query. We tested a couple of cosine and sigmoid functions for computing vector similarity of the projected vector and the feedback terms. The experimental results reveal that using the cosine similarity and a softmax layer works as well as the state-of-the-art feedback techniques and even better in a few number of collections. ECDMM has significant improvements up to 3.8% compared to the state-of-the-art models in MAP.

This work inspires a number of future works. First, we want to study the usage of the proposed method in low-dimensional profile modelling in recommendation systems. In recommender systems, there is stream of documents being proposed to the users and thus we can update the low-dimensional profile vectors incrementally. The proposed formulation of low-dimensional query updating can be adapted to this work. Although the main focus of this work is to provide a robust formulation for query modeling/updating, tailoring semantic networks (e.g., WordNet, Concept-Net.) for positive sampling seems to be interesting. Therefore, our second goal is to study the effect of negative sampling on difficult queries.

6. ACKNOWLEDGMENT

The authors would like to thank Hamed Zamany from Google corporation for his helpful comments on this work.

7. REFERENCES

- [1] M. AlMasri, C. Berrut, and J.-P. Chevallet. *A Comparison of Deep Learning Based Query Expansion with Pseudo-Relevance Feedback and Mutual Information*, pages 709–715. Springer International Publishing, Cham, 2016.
- [2] Y. Bengio and G. Corrado. Bilbowa: Fast bilingual distributed representations without word alignments. 2014.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *JMLR*, 3(Jan):993–1022, 2003.
- [4] S. Clinchant and F. Perronnin. Aggregating continuous word embeddings for information retrieval. In *Proceedings of the Workshop on Continuous Vector Space Models and their Compositionality*, pages 100–109, 2013.
- [5] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing (almost) from scratch. *JMLR*, 12:2493–2537, 2011.
- [6] J. Dadashkarimi, M. S. Shahshahani, A. Tebbifakhr, H. Faili, and A. Shakery. Dimension projection among languages based on pseudo-relevant documents for query translation. *arXiv preprint arXiv:1605.07844*, 2016.
- [7] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *JASIST*, 41(6):391, 1990.
- [8] D. Ganguly, D. Roy, M. Mitra, and G. J. Jones. Word embedding based generalized language model for information retrieval. In *SIGIR '15*, SIGIR '15, pages 795–798, New York, NY, USA, 2015. ACM.
- [9] Y. Goldberg and O. Levy. word2vec explained: Deriving mikolov et al.’s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722*, 2014.
- [10] M. Grbovic, N. Djuric, V. Radosavljevic, and N. Bhamidipati. Search retargeting using directed query embeddings. In *WWW '15 Companion*, pages 37–38, New York, NY, USA, 2015. ACM.
- [11] H. B. Hashemi and A. Shakery. Mining a Persian-English Comparable Corpus for Cross-language Information Retrieval. *IP&M*, 50(2):384–398, 2014.
- [12] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [13] T. Hofmann. Probabilistic latent semantic indexing. In *SIGIR '99*, pages 50–57. ACM, 1999.
- [14] N. A. Jaleel, J. Allan, W. B. Croft, F. Diaz, L. S. Larkey, X. Li, M. D. Smucker, and C. Wade. Umass at trec 2004: Novelty and hard. In *TREC*, 2004.
- [15] R. Kiros, R. S. Zemel, and R. Salakhutdinov. A multiplicative model for learning distributed text-based attribute representations. *CoRR*, abs/1406.2710, 2014.
- [16] M. Kusner, Y. Sun, N. Kolkin, and K. Q. Weinberger. From word embeddings to document distances. In D. Blei and F. Bach, editors, *ICML '15*, pages 957–966. JMLR Workshop and Conference Proceedings, 2015.
- [17] V. Lavrenko and W. B. Croft. Relevance-based language models. In *SIGIR '01*, pages 120–127, 2001.
- [18] Q. V. Le and T. Mikolov. Distributed representations of sentences and documents. *arXiv preprint arXiv:1405.4053*, 2014.
- [19] Q. V. Le and T. Mikolov. Distributed representations of sentences and documents. *CoRR*, abs/1405.4053, 2014.
- [20] Y. Lv and C. Zhai. Revisiting the divergence minimization feedback model. In *CIKM '14*, pages 1863–1866, 2014.
- [21] C. D. Manning. Understanding human language: Can NLP and deep learning help? In *SIGIR '16*, page 1, 2016.
- [22] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Adv. in Neur. Inf. Proc. Sys.*, pages 3111–3119, 2013.
- [23] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *EMNLP '14*, pages 1532–1543, 2014.
- [24] N. Rahimi, A. Shakery, J. Dadashkarimi, M. Aryannejad, M. Dehghani, and H. N. Esfahani. Building a multi-domain comparable corpus using a learning to rank method. *NLE*, 22(4):627–653, 2016.
- [25] R. Socher, J. Pennington, E. H. Huang, A. Y. Ng, and C. D. Manning. Semi-supervised recursive autoencoders for predicting sentiment distributions. In *EMNLP '11*, pages 151–161, 2011.
- [26] R. Socher, A. Perelygin, J. Y. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP '13*, volume 1631, page 1642, 2013.
- [27] K. Sparck Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation*, 28(1):11–21, 1972.
- [28] J. Turian, L. Ratinov, and Y. Bengio. Word representations: a simple and general method for semi-supervised learning. In *ACL '10*, pages 384–394. ACL, 2010.
- [29] I. Vulic and M. Moens. Monolingual and cross-lingual information retrieval models based on (bilingual) word embeddings. In *SIGIR '15*, pages 363–372, 2015.
- [30] X. Wang, H. Fang, and C. Zhai. A study of methods for negative relevance feedback. In *SIGIR '08*, pages 219–226, New York, NY, USA, 2008. ACM.
- [31] H. Zamani, J. Dadashkarimi, A. Shakery, and W. B. Croft. Pseudo-relevance feedback based on matrix factorization. In *CIKM '15*, Indiana Police, USA, 2015. ACM.
- [32] C. Zhai and J. Lafferty. Model-based Feedback in the Language Modeling Approach to Information Retrieval. In *CIKM '01*, pages 403–410, Atlanta, Georgia, USA, 2001.
- [33] G. Zheng and J. Callan. Learning to reweight terms with distributed representations. In *SIGIR '15*, pages 575–584, New York, NY, USA, 2015. ACM.